# Beyond Keyword Search: Discovering Relevant Scientific Literature

Khalid El-Arini     Carlos Guestrin

**ML**

**MACHINE LEARNING**
**D E P A R T M E N T**

**Carnegie Mellon**®

# Report Documentation Page

| 1. REPORT DATE **JUN 2011** | 2. REPORT TYPE | 3. DATES COVERED **00-00-2011 to 00-00-2011** |
|---|---|---|

| 4. TITLE AND SUBTITLE **Beyond Keyword Search: Discovering Relevant Scientific Literature** | | 5a. CONTRACT NUMBER |
|---|---|---|
| | | 5b. GRANT NUMBER |
| | | 5c. PROGRAM ELEMENT NUMBER |
| 6. AUTHOR(S) | | 5d. PROJECT NUMBER |
| | | 5e. TASK NUMBER |
| | | 5f. WORK UNIT NUMBER |

| 7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) **Carnegie Mellon University,School of Computer Science,Machine Learning Department,Pittsburgh,PA,15213** | 8. PERFORMING ORGANIZATION REPORT NUMBER |
|---|---|

| 9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) | 10. SPONSOR/MONITOR'S ACRONYM(S) |
|---|---|
| | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) |

12. DISTRIBUTION/AVAILABILITY STATEMENT
**Approved for public release; distribution unlimited**

13. SUPPLEMENTARY NOTES

14. ABSTRACT
**In scientific research, it is often difficult to express information needs as simple keyword queries. We present a more natural way of searching for relevant scientific literature. Rather than a string of keywords, we define a query as a small set of papers deemed relevant to the research task at hand. By optimizing an objective function based on a fine-grained notion of influence between documents, our approach efficiently selects a set of highly relevant articles. Moreover, as scientists trust some authors more than others, results are personalized to individual preferences. In a user study, researchers found the papers recommended by our method to be more useful trustworthy and diverse than those selected by popular alternatives, such as Google Scholar and a state-of-the-art topic modeling approach.**

15. SUBJECT TERMS

| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES | 19a. NAME OF RESPONSIBLE PERSON |
|---|---|---|---|---|---|
| a. REPORT **unclassified** | b. ABSTRACT **unclassified** | c. THIS PAGE **unclassified** | **Same as Report (SAR)** | **35** | |

# Beyond Keyword Search: Discovering Relevant Scientific Literature

**Khalid El-Arini**        **Carlos Guestrin**

June 2011

CMU-ML-11-102

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

## Abstract

In scientific research, it is often difficult to express information needs as simple keyword queries. We present a more natural way of searching for relevant scientific literature. Rather than a string of keywords, we define a query as a small set of papers deemed relevant to the research task at hand. By optimizing an objective function based on a fine-grained notion of influence between documents, our approach efficiently selects a set of highly relevant articles. Moreover, as scientists trust some authors more than others, results are personalized to individual preferences. In a user study, researchers found the papers recommended by our method to be more useful, trustworthy and diverse than those selected by popular alternatives, such as Google Scholar and a state-of-the-art topic modeling approach.

# 1    Introduction

For generations, scientists have built upon the published work of their predecessors and contemporaries in order to make new discoveries. However, as the number of publications has grown, it has become increasingly difficult for scientists to find relevant prior work for their particular research. In fact, as early as 1755, the French philosopher Denis Diderot presciently forewarned that there would come a day when "it will be almost as convenient to search for some bit of truth concealed in nature as it will be to find it hidden away in an immense multitude of bound volumes," [14]. Today, we can quantify this "immense multitude" to include tens of millions of articles published in tens of thousands of journals and conferences [44].

Currently, researchers primarily rely on keyword search of online indices such as Google Scholar and PubMed to help them combat this overload of information. While these tools are indispensable, there are many instances where a researcher's information need cannot be easily specified as a simple string of keywords. Often, such a keyword query is either overly broad, returning many articles that are at best loosely related to the researcher's specific need, or too narrow, potentially returning no articles at all. In these occasions, it may be more natural for the scientist to specify his query as a small set of papers rather than as a set of words. In particular, having already read some articles that are related to the specific task at hand, the scientist can ask, "given that these papers represent my immediate research focus, what else should I read?".

Here, we present an algorithm for discovering relevant scientific literature by responding to queries of this form. More formally, given a small set of papers $\mathcal{Q}$ that we refer to as the *query set*, we seek to return a set $\mathcal{A}$ of additional papers that are related to the concept defined by the query. Intuitively, a paper that cites all of the articles in $\mathcal{Q}$ is likely to represent related research. Likewise, a paper that is cited by every article in $\mathcal{Q}$ might contain relevant background information. However, it is restrictive to require the papers in $\mathcal{A}$ to have a direct citation to or from every article in the query set, as such papers are not guaranteed to exist. Instead, we wish to select a set $\mathcal{A}$ that maximizes a more general notion of *influence* to and from the papers in $\mathcal{Q}$.

# 2    Modeling Scientific Influence

To define a notion of influence in scientific literature, we observe that the content of a publication is an amalgam of several sources, combining cited prior work with the authors' novel insights and background experience. For a given collection of articles, ideas travel *from cited papers to citing papers*, and from earlier to subsequent papers by the same author (Figure 1A). Our notion of influence should capture this transfer of ideas, modeling both the extent to which ideas travel between documents, as well as their topical matter. To achieve such fine-grained detail, we define influence with respect to the *individual concepts* found in a document collection, which could be, e.g., technical terms or informative phrases.[1] For example, we might say that the ideas transferred from one paper to another involve the concepts "energy" or "nitric oxide."

For each concept $c$ in our vocabulary of concepts $\mathcal{C}$, we define a directed, acyclic graph $G_c$, where the nodes represent papers that contain $c$ and the edges represent citations and common authorship. Figures 1B and 1C show two such graphs for a subset of articles from the Proceedings of the National Academy of Sciences (PNAS), for the concepts "plant" and "stress." While a path between two nodes in such a graph may indicate influence with respect to a particular concept, mere existence of a path does little to express the *degree* to which this influence occurs. To capture the degree of influence, we define a weight $\theta^{(c)}_{x \to y}$ for each edge $(x, y)$ in graph $G_c$, representing the probability of *direct* influence from paper $x$ to paper $y$ with respect to concept $c$. We can then use these edge weights to define a probabilistic, concept-specific notion of influence between any two papers in the document collection.
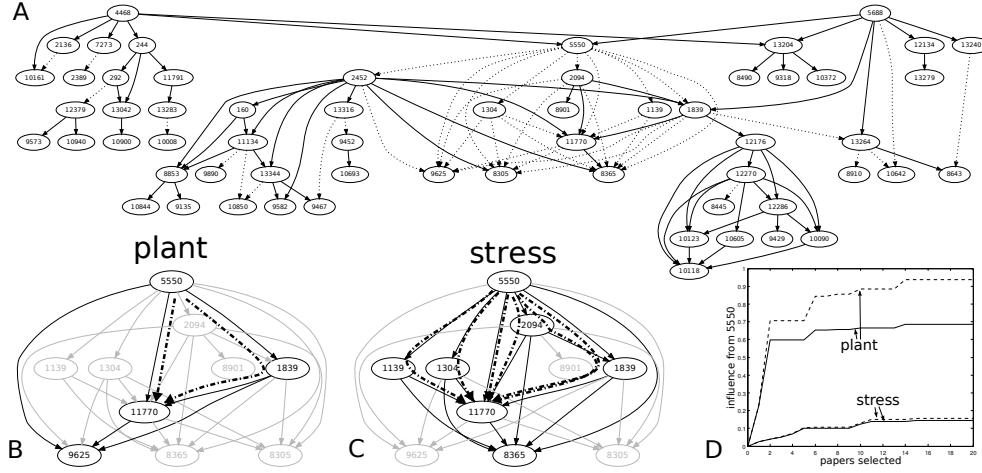
Figure 1: **(A)** A graph of articles from the Proceedings of the National Academy of Sciences (PNAS). Nodes represent papers, solid edges represent citations ($x \to y$ if $y$ cites $x$) and dotted edges represent common authorship ($x \to y$ if $x$ is older than $y$ and $x, y$ share an author). More details on the data sets used in this paper can be found in the appendix. **(B,C)** Subgraphs of (A), limited to papers containing the concepts "plant" and "stress," respectively (other papers are grayed out). Thick dashed lines indicate paths of influence between papers 5550 and 11770. **(D)** Example illustrating how Equation 1 penalizes redundancy. The first two papers selected exhibit a high influence with respect to "plant," and thus subsequently adding such papers to $\mathcal{A}$ causes little increase in Equation 1 (solid lines), especially when compared to the sum of individual influences (dashed lines). The influence with respect to "stress" remains low, thus never triggering such a redundancy penalty.

## 2.1 Defining edge weights

Figure 2 shows an example from the PNAS data set illustrating how we define the weight $\theta^{(c)}_{x \to y}$ on each edge. Here, article 9467 cites two articles containing the concept "oxygen," $\{424, 13344\}$, indicated by the solid black edges. The dotted black edges indicate that two other articles, $\{1829, 7657\}$, contain the concept "oxygen" and share authors with 9467. (The dotted gray edge indicates that there is a third article sharing authors with 9467 that *does not* contain "oxygen.") We assume that every occurrence of the concept "oxygen" in 9467 is either a novel idea or is directly inspired by one of these sources. Thus, we view the weight $\theta^{(c)}_{x \to y}$ as the probability a random instance of concept $c$ in paper $y$ was directly inspired by paper $x$.

The bar graph over the nodes in Figure 2 illustrates the proportion of the content of each paper consisting of the "oxygen" concept. For instance, the height of the first bar on the left is $n^{(oxygen)}_{424}/N_{424}$, where $n^{(c)}_x$ is the frequency of concept $c$ in document $x$, and $N_x = \sum_{c \in \mathcal{C}} n^{(c)}_x$ is the total length of document $x$. Additionally, the bars over 1829 and 7657 are shortened to one third of their original height (indicated in light gray), representing the intuition that an explicit citation is a more informative relationship than common authorship. The authors of 9467 have three prior publications in this example, and thus by dividing by three, the effective total contribution of these papers is that of a single paper. Finally, we represent the novelty distribution for a particular paper $y$ as the average distribution over concepts for all papers published in the same year as $y$. In this case, the novelty contribution for "oxygen" is dominated by the four papers. (We note that there are no actual novelty nodes in the graph, as the associated distribution is only used for normalization.)

Here, $\theta^{(oxygen)}_{x \to 9467}$ is proportional to the height of the corresponding bar in the plot. More generally, if a paper $y$ cites papers $\{r_1, \dots, r_k\}$, and the authors have previously written papers $\{b_1, \dots, b_l\}$, then the edge weights are

---

[1]For our experiments, we use a simple tf-idf heuristic to extract informative words which we use as concepts, as described in the appendix.
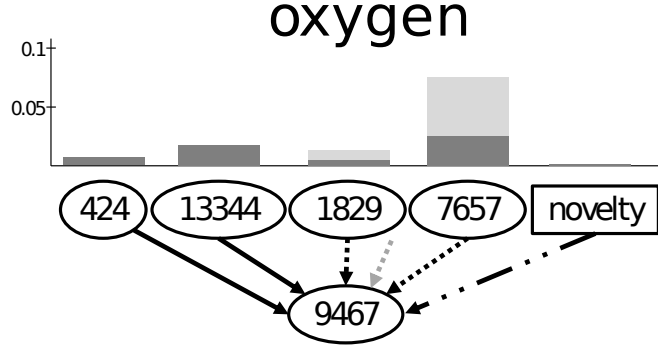
Figure 2: An example from the PNAS data set, illustrating the edge weight computation for a node in $G_{oxygen}$. Solid black edges indicate citations, while dotted black edges indicate common authorship. The dotted gray edge refers to a paper sharing an author with 9467, but not containing the concept "oxygen." Edge weights are assigned proportional to the bar chart, indicating the prevalence of "oxygen" in each parent node. The bars over 1829 and 7657 are shortened to one third of their original height (indicated in light gray), such that the contribution due to common authorship is equivalent to that of a single paper. The novelty node is only used to normalize the edge weights, and in this case is dominated in influence by the other articles.

defined as follows:

$$
\theta_{r_i \to y}^{(c)} = \frac{1}{Z} \frac{n_{r_i}^{(c)}}{N_{r_i}},
$$

$$
\theta_{b_i \to y}^{(c)} = \frac{1}{Z} \frac{n_{b_i}^{(c)}}{l \cdot N_{b_i}},
$$

with normalization constant,

$$
Z = \sum_{j=1}^{k} \frac{n_{r_j}^{(c)}}{N_{r_j}} + \frac{1}{l} \sum_{j=1}^{l} \frac{n_{b_j}^{(c)}}{N_{b_j}} + novel_y^{(c)},
$$

where $novel_y^{(c)}$ is the average proportion of concept $c$ across all papers published in the same year as $y$.

## 2.2 Calculating influence

Given a concept-specific weight for each edge in the citation graph, representing the *direct* influence between two neighboring nodes, we can now define the influence between any two papers in our collection. In particular, if we say that each edge $x \to y$ in $G_c$ is *active* with some probability $\theta_{x \to y}^{(c)}$, we arrive at the following definition:

**Definition 1.** *The influence between papers $u$ and $v$ with respect to concept $c$, $Influence_c(u \leftrightarrow v)$, is the probability there exists a directed path in $G_c$ from one paper to the other, consisting only of active edges.*

While intuitive, the exact computation of this probability is intractable, as the problem of computing connectedness in a random graph belongs to the #P-complete class of computational problems [46, 38], for which there are no known polynomial-time solutions. We can overcome this computational hurdle via approximation, by employing one of two methods: 1) a simple Monte Carlo sampling procedure with theoretical guarantees; and, 2) a deterministic, linear-time dynamic programming heuristic, based on the assumption that the paths between two nodes are independent of each other.

### 2.2.1 Sampling

The simplest procedure for estimating the influence between two nodes is to generate samples directly based on the definition of influence. Each sample is generated as follows:

For each concept $c$:

1. Mark each edge $x \to y$ in $G_c$ as active with probability $\theta_{x \to y}^{(c)}$.
2. For all pairs of nodes $(u, v)$, record whether a path exists between them using only active edges.

After generating $B$ samples, the probability that a node $u$ influences a node $v$ with respect to concept $c$ is simply estimated as the proportion of the $B$ samples for concept $c$ in which an active path from $u$ to $v$ exists. A natural question to ask is, how many samples do we need for a reasonable estimate of influence? A short proof using Hoeffding's Inequality shows us that the number of samples we need grows only *logarithmically* with the number of articles in the document collection.

**Theorem 1.** *In order to estimate $m$ influence values such that, with probability $\eta$, each of the $m$ estimates is no more than $\delta$ away from its true value, a sufficient number of samples $B$ is $\frac{2}{\delta^2} \log(2m/\delta)$.*

*Proof.* We wish to estimate $m$ influence probabilities, $p_1, p_2, \ldots, p_m$, using $\hat{p}_1, \hat{p}_2, \ldots, \hat{p}_m$, where $\hat{p}_j = \frac{1}{B} \sum_{i=1}^{B} X_j^{(i)}$, and each $X_j^{(i)}$ is a random variable that is either 0 or 1, representing whether the $j$th pair of nodes is connected via an active path in sample $i$. Note that by our definition of influence, $E[\hat{p}_j] = \frac{1}{B} \sum_{i=1}^{B} E[X_j^{(i)}] = p_j$.

We let $\epsilon_j = |\hat{p}_j - p_j|$, the absolute difference between influence value $p_j$ and its estimate using the sampling methodology from above. Given some $\delta$, we want $P(\epsilon_1 \geq \delta \vee \epsilon_2 \geq \delta \vee \ldots \vee \epsilon_m \geq \delta)$ to be small.

$$
\begin{aligned}
P\left( \bigvee_{j=1}^{m} (\epsilon_j \geq \delta) \right) &\leq \sum_{j=1}^{m} P\left( \epsilon_j \geq \delta \right) \\
&= \sum_{j=1}^{m} P\left( |\hat{p}_j - p_j| \geq \delta \right) \\
&= \sum_{j=1}^{m} P\left( \left| \frac{1}{B} \sum_{i=1}^{B} X_j^{(i)} - \frac{1}{B} \sum_{i=1}^{B} E\left[ X_j^{(i)} \right] \right| \geq \delta \right) \\
&= \sum_{j=1}^{m} P\left( \left| \frac{1}{B} \sum_{i=1}^{B} X_j^{(i)} - E\left[ X_j^{(i)} \right] \right| \geq \delta \right) \\
&= \sum_{j=1}^{m} P\left( \left| \sum_{i=1}^{B} X_j^{(i)} - B \cdot E\left[ X_j^{(i)} \right] \right| \geq B\delta \right) \\
&\leq \sum_{j=1}^{m} 2 \exp\left( \frac{-2B^2 \delta^2}{4B} \right) \\
&= 2m \exp\left( \frac{-B\delta^2}{2} \right),
\end{aligned}
$$

where the first inequality is due to the union bound, and the second inequality is due to Hoeffding.

Thus, the probability that any of our $m$ estimates is more than $\delta$ away from its true value given $B$ samples is

less than or equal to $2m \exp(-B\delta^2/2)$. For this probability to be less than or equal to $\eta$, we need:

$$
\begin{aligned}
2m \exp\left(\frac{-B\delta^2}{2}\right) &\leq& \eta, \\
\exp\left(\frac{-B\delta^2}{2}\right) &\leq& \frac{\eta}{2m}, \\
-B\delta^2 &\leq& 2\log\left(\frac{\eta}{2m}\right), \\
B &\geq& \frac{-2}{\delta^2}\log\left(\frac{\eta}{2m}\right) \\
&=& \frac{2}{\delta^2}\log\left(\frac{2m}{\eta}\right).
\end{aligned}
$$

$\square$

As the number of influence values to estimate is quadratic in the number of articles, the number of samples we need is logarithmic in the total number of articles. While this is a heartening result, we find that for large document collections, generating enough samples can still be a time-consuming process.

### 2.2.2 Independence heuristic

As an alternative to sampling, we describe an efficient dynamic programming heuristic based on the assumption that the paths between two nodes in $G_c$ are independent of each other. For instance, in Figure 1B, the two influence paths between 5550 and 11770 with respect to the concept "plant" are completely independent of each other. Thus, the probability of at least one active path existing between the two nodes in this situation can be computed exactly:

$$
\begin{aligned}
&Influence_{plant}(5550 \to 11770) \\
&= 1 - P(\text{there is no influence between 5550 and 11770}) \\
&= 1 - P(\text{there is no direct influence from 5550}) \cdot \\
&\quad\ P(\text{there is no influence through 1839}) \\
&= 1 - (1 - \theta^{(plant)}_{5550 \to 11770})(1 - \theta^{(plant)}_{5550 \to 1839}\theta^{(plant)}_{1839 \to 11770}).
\end{aligned}
$$

The second equality follows from the independence of the two paths. On the other hand, looking at Figure 1C, we find the paths between the two nodes in $G_{stress}$ are not independent, making such a calculation more problematic.

Based on this intuition, if we rashly assume that the paths between two nodes will *always* be independent of each other in $G_c$, for all $c$, we arrive at a simple, efficient heuristic for computing the influence between all pairs of nodes (Algorithm 1). By traversing the graph in topological order, we know that when we arrive at a node we will have already computed all the influence going to its parents. Using these influences and our independence assumption, we can then immediately compute the influence to the node itself. We note that this algorithm requires the graphs to be acyclic.[2]

While the independence assumption upon which this heuristic is based certainly is not true in general, we find that, nevertheless, the values we compute are close to what we would expect from sampling (cf. Figure 3). Thus, despite not being amenable to theoretical guarantees, we find this heuristic works well in practice.

---

[2]Based on simple chronology, one would expect a citation graph to be acyclic; after all, a researcher cannot cite a paper if it does not yet exist. However, this is not quite the case in practice (e.g., colleagues writing papers simultaneously may cite each other). Details on how we address this problem can be found in the appendix.

---
**Algorithm 1** Dynamic Programming Heuristic for Influence

---
$N$: number of documents
$\mathcal{C}$: vocabulary of concepts
// Initialize to empty 3D array
// $influenceEstimate[c][x][y]$ will contain influence
//    from $x$ to $y$ with respect to concept $c$.
$influenceEstimate \leftarrow array[|\mathcal{C}|][N][N]$
**for all** $c \in \mathcal{C}$ **do**
    **for all** nodes $y$ in $G_c$ **do**
      // Initialize to identity
      $influenceEstimate[c][y][y] \leftarrow 1$
    $topoOrder \leftarrow$ topological order of nodes in $G_c$
    **for** $y \in topoOrder$ **do**
      // $influenceEstimate[c][][x]$ already calculated
      //    for all $x \in parents(y)$
      **if** $parents(y) = \emptyset$ **then**
        continue
      $influenceFromParents \leftarrow array[|parents(y)|]$
      **for all** $x \in parents(y)$ **do**
        // Influence to the parent multiplied by
        //    the edge weight
        $influenceFromParents[x] \leftarrow$
          $influenceEstimate[c][][x] \cdot \theta_{x \to y}^{(c)}$
      // Product is element-wise
      $influenceEstimate[c][][y] \leftarrow$
        $1 - \prod_{x \in parents(y)}(1 - influenceFromParents[x])$

---

## 2.3 Selecting papers

As motivated in Section 1, given a query set of papers $\mathcal{Q}$, we wish to select a small set of related papers $\mathcal{A}$ that exhibit a high degree of *influence* to or from the query set. Moreover, the set of papers we select should be both *relevant* and *diverse*.

### 2.3.1 Relevance

The influence between the query set $\mathcal{Q}$ and the result set $\mathcal{A}$ should be focused on the concepts that are important or prevalent in both sets of documents. First, to ensure that the selected documents pertain to the concepts most prevalent in the query set, we define a weight $\gamma_q^{(c)}$ proportional to the frequency of concept $c$ in query document $q$.

Likewise, from the perspective of the result set, a document $d$ might contain a single occurrence of the concept "plant," and that single occurrence might be heavily influenced by one of the query documents $q$. However, as $d$ only tangentially mentions "plant," we do not wish this strong influence to incentivize its inclusion in the result set. Thus, we define a probability $\beta_d^{(c)}$ indicating the importance of concept $c$ in document $d$. Specifically, we define this as the probability a concept $c$ is observed in a finite number $\ell$ of independent samples (with replacement) from the document's word distribution: $\beta_d^{(c)} = 1 - (1 - \gamma_d^{(c)})^\ell$. (Here, $\ell$ is a parameter of our model that we set to 20 in our experiments.)

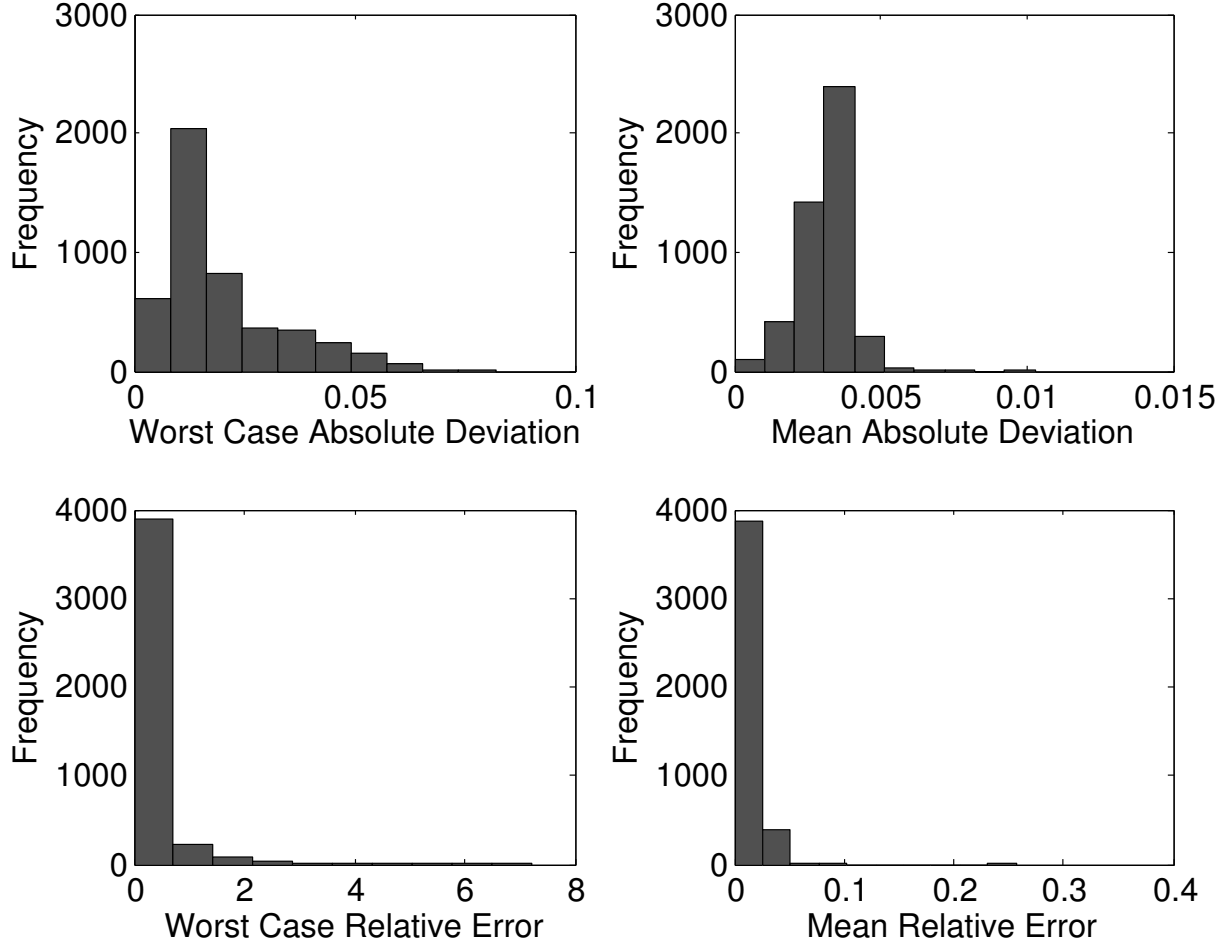Figure 4 provides an illustrative example of these weights.

Figure 3: This figure shows a comparison on the PNAS data set between the influence values computed via sampling ($B = 6530$) and those computed using the independence heuristic. For all concepts and all pairs of articles with meaningful influence between them (i.e., not trivially zero, as is the case when the nodes are not connected in the graph), we compute the influence using both methods, and record the absolute deviation ($|sampling - heuristic|$) and relative error ($|sampling - heuristic|/sampling$). The worst case and mean values of these measures for each concept are plotted above. For this setting of $B$, the estimates computed via sampling are likely ($> 95\%$) to be within 0.075 of their true values.

### 2.3.2 Diversity

Diversity is important in this setting as it is difficult to predict the exact information need of a researcher, and thus providing a wide variety of papers increases the likelihood of query satisfaction. As such, we define the influence between a single query paper $q \in \mathcal{Q}$ and a set of documents $\mathcal{A}$ in a manner that penalizes redundancy in the result set, thereby promoting diversity. Specifically, if we define this *set influence*, $Influence_c(q \leftrightarrow \mathcal{A})$, as the probability influence exists between $q$ and *at least one* document in $\mathcal{A}$, we create a disincentive for $\mathcal{A}$ to contain multiple papers with similar influence patterns to and from $q$; such a redundant set $\mathcal{A}$ would exhibit less influence
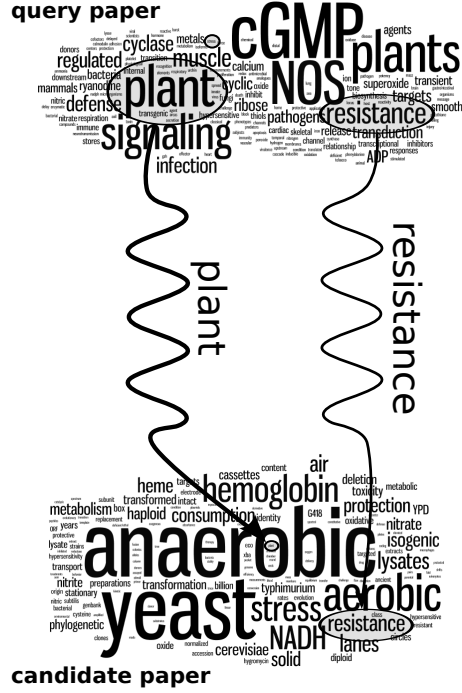
Figure 4: The top cloud represents a query paper (5550), the bottom word cloud represents a paper to be selected (11770) and the lines between them represent individual influences of varying strength. In each word cloud, the size of a word is proportional to its frequency in the corresponding article. $\gamma$ is illustrated by the shaded ellipses in the top word cloud, showing a higher incentive to pick articles about "plant" or "resistance" than about "stress." However, despite its prevalence in the query document, "plant" is only tangentially present in article 11770, and thus $\beta$ ensures a low degree of influence. This can be contrasted with "resistance," which is prevalent in both documents and displays a high degree of influence.

than one composed of a broader set of documents. Formally,

$$
\begin{aligned}
Influence_c(q \leftrightarrow \mathcal{A}) = \\
1 - \prod_{d \in \mathcal{A}} \left( 1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)} \right).
\end{aligned}
\tag{1}
$$

We note the use of the probability $\beta_d^{(c)}$ here to safeguard against selecting documents that are only tangentially related to the important concepts in the query papers.

Figure 1D shows an example illustrating how the marginal gain in set influence with respect to the concept "plant" diminishes as more papers are added to the result set $\mathcal{A}$. In particular, beyond a certain level of influence, the gain observed in Equation 1 from adding additional documents to the result set is smaller than would be expected if we were naïvely summing the individual influences. We do not see the same redundancy penalty with respect to "stress," as the result set is not sufficiently influenced with respect to this concept.

### 2.3.3 Optimization

Given this definition of set influence, we can now define an objective function that, when maximized, returns a diverse set of papers highly relevant to the query:

$$F_{\mathcal{Q}}(\mathcal{A}) \quad = \quad \sum_{q \in \mathcal{Q}} \sum_{c \in \mathcal{C}} \gamma_q^{(c)} Influence_c(q \leftrightarrow \mathcal{A}). \tag{2}$$

While, in general, solving such a combinatorial optimization problem is intractable, Equation 2 exhibits an intuitive diminishing returns property known as *submodularity*, allowing for efficient near-optimal solutions.

**Definition 2** (Submodularity). *A set function F is* submodular *if,* $\forall \mathcal{A} \subseteq \mathcal{B} \subseteq \mathcal{V}, \forall s \in \mathcal{V} \setminus \mathcal{B}, F(\mathcal{A} \cup \{s\}) - F(\mathcal{A}) \geq F(\mathcal{B} \cup \{s\}) - F(\mathcal{B}).$

Intuitively, this means that the utility of adding a particular paper to a result set decreases as the result set gets larger.

**Theorem 2.** *Equation 2 is submodular and monotonic.*

*Proof.* Let $\mathcal{B} \subseteq \mathcal{V}$ and $s \in \mathcal{V} \setminus \mathcal{B}$.

$$Influence_c(q \leftrightarrow \mathcal{B} \cup \{s\}) - Influence_c(q \leftrightarrow \mathcal{B})$$

$$= 1 - \prod_{d \in \mathcal{B} \cup \{s\}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) - \left(1 - \prod_{d \in \mathcal{B}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)})\right)$$

$$= \prod_{d \in \mathcal{B}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) - \prod_{d \in \mathcal{B} \cup \{s\}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)})$$

$$= \prod_{d \in \mathcal{B}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) \left(1 - (1 - Influence_c(q \leftrightarrow s)\beta_s^{(c)})\right)$$

$$= \prod_{d \in \mathcal{B}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) \left(Influence_c(q \leftrightarrow s)\beta_s^{(c)}\right).$$

Because $Influence_c(q \leftrightarrow d)$ and $\beta_d^{(c)}$ are defined as probabilities, and thus lie in the range $[0, 1]$, we know that this quantity is non-negative, making Equation 2 monotonic. Moreover, for any $\mathcal{A} \subseteq \mathcal{B}$, we have that,

$$\prod_{d \in \mathcal{B}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) \leq \prod_{d \in \mathcal{A}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}).$$

Hence,

$$\prod_{d \in \mathcal{B}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) \left(Influence_c(q \leftrightarrow s)\beta_s^{(c)}\right)$$

$$\leq \prod_{d \in \mathcal{A}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) \left(Influence_c(q \leftrightarrow s)\beta_s^{(c)}\right)$$

$$= \prod_{d \in \mathcal{A}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)}) - \prod_{d \in \mathcal{A} \cup \{s\}} (1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)})$$

$$= Influence_c(q \leftrightarrow \mathcal{A} \cup \{s\}) - Influence_c(q \leftrightarrow \mathcal{A}).$$

Thus, $Influence_c(q \leftrightarrow \mathcal{A})$ is submodular. Since submodularity is closed under non-negative linear combinations, and our weights $\gamma_q^{(c)} \geq 0$, it directly follows that our objective function in Equation 2 is submodular. $\square$
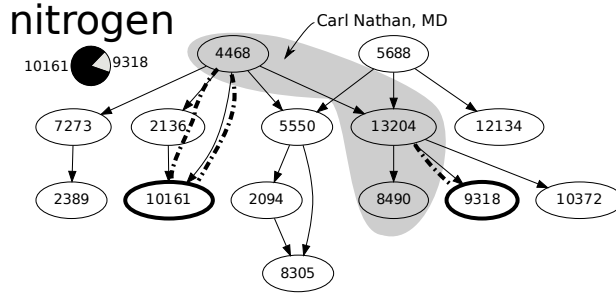
Figure 5: Example illustrating trust calculation for an immunologist asking, "How much do I trust Carl Nathan with respect to the concept 'nitrogen'?" Thick dashed lines indicate influence from Dr. Nathan to individual elements of $\mathcal{B}$, and pie chart represents relative prevalence of the word "nitrogen" in the two papers in $\mathcal{B}$.

Although maximizing submodular functions is NP-hard [26], by discovering this property in our problem, we can take advantage of several efficient approximation algorithms with theoretical guarantees. For example, the classic result of Nemhauser et al. [32] shows that by simply applying a greedy algorithm to maximize our objective function, we can obtain a $(1 - \frac{1}{e})$ approximation of the optimal value. Thus, a simple greedy optimization can provide us with a near-optimal solution. However, since our set of articles is very large, a naïve greedy approach can be too costly. Therefore, we use CELF [29], which provides the same approximation guarantees, but uses lazy evaluations, often leading to dramatic speedups.

## 3   Trust and Personalization

Considering our running example of PNAS articles (Figure 1A), we can set our query set to be $\mathcal{Q} = \{4468, 5688\}$, the parents of "Nitric Oxide in Plant Immunity" (5550). Optimizing Equation 2 for this query produces a result set of articles ranging in topics from plant biology to immunology (cf. Table 2). While these articles may be relevant to the query, a major shortcoming is that every researcher who submits this query will receive an identical result set. For any given topic, different researchers trust different authors and publications, and the objective in Equation 2 provides no means to express these preferences. While a long line of prior work exists on summarizing the impact of an author or publication with a single number [2], often based on citation statistics [20, 25] or eigenvector methods [27, 36, 12, 39], here we wish to capture a more detailed picture of the relationship between a researcher and the authors he cites.

In order to properly model such an individual notion of *trust* in the setting of scholarly research, we consider two motivating scenarios:

1. Different authors command different levels of respect from their research communities, e.g., a Nobel laureate versus a first-year graduate student, as an extreme case.
2. Even among distinguished scientists, a particular researcher's interests may be aligned more closely with some than others. Thus, beyond simply differentiating novices from experts, a notion of trust should also capture differences in research interests. For example, asking computer scientists to name whom they most associate with the concept "network" may yield Judea Pearl (Bayesian networks), Jon Kleinberg (social networks), Geoff Hinton (neural networks) or Van Jacobson (computer networks), depending on who is answering. All are distinguished researchers, but each is associated with a distinct subfield of computer science.

At the heart of both scenarios is a personal question that is often answered differently by different researchers: *How much do I trust this author with respect to this concept?*

By answering this question, a researcher would enable us to formally incorporate his trust preferences into our objective function, allowing us to select papers tailored specifically to his tastes. However, as researchers will not be able to provide an answer for every combination of authors and concepts, we must elicit their trust preferences in a less onerous manner. In order to do so, we assume that trust is *transitive*. For example, if Alice trusts an article, and that article is heavily influenced by Bob with respect to the concept "network," then Alice is likely to also trust Bob with respect to "network." Thus, at a fundamental level, a researcher need only specify a set of trusted papers $\mathcal{B}$, from which we can infer answers to the above question. As a shortcut, a researcher may choose to define $\mathcal{B}$ indirectly by specifying a list of trusted journals and conferences, or subsets thereof (e.g., a particular conference track or article classification). $\mathcal{B}$ could also be specified as the papers *cited by* one or more trusted authors, representing a look at one's research through the eyes of another scientist, potentially in another field. Thus, a plucky physicist could ask, "What would Steven Chu recommend I read?", and obtain a set of papers related to his query, yet tailored to the research interests and trust preferences of the Nobel laureate.

With this intuition in mind, we define $\tau_{a|\mathcal{B}}^{(c)}$, the probability a researcher trusts author $a$ with respect to concept $c$, given trusted articles $\mathcal{B}$. (The "$|\mathcal{B}$" notation in this section indicates personalizing with respect to trusted set $\mathcal{B}$.) Figure 5 illustrates how we compute $\tau_{a|\mathcal{B}}^{(c)}$ for a particular example from PNAS, where the concept $c$ is "nitrogen," the author $a$ is Carl Nathan, MD, and the researcher has specified two immunology papers as his trusted set, $\mathcal{B} = \{10161, 9318\}$. For each paper $b \in \mathcal{B}$, we compute how much the author $a$ influenced $b$ with respect to concept $c$. As our influence is now expressed from an *author* to an article, we treat all of an author's papers as a single unit.

**Definition 3.** *The influence from author $a$ to article $b$ with respect to concept $c$, $AuthorInfluence_c(a \to b)$, is the probability there exists a directed path in $G_c$ from* any *article written by $a$ to article $b$ consisting only of active edges, where each edge is (independently) active with probability $\theta_{x \to y}^{(c)}$.*

As before, we employ sampling or dynamic programming to efficiently estimate this otherwise intractable computation (cf. Algorithm 2).

In our example, we first look at how much Dr. Nathan's papers influence 10161 with respect to "nitrogen," and again from Dr. Nathan's papers to 9318. We now weigh these two influences by the prevalence of the word "nitrogen" in each paper $b$ (as indicated by the pie chart in Figure 5), and define $\tau_{a|\mathcal{B}}^{(c)}$ to be the weighted sum of the two.

More generally, we have:

$$\tau_{a|\mathcal{B}}^{(c)} = \begin{cases} \frac{1}{N_{\mathcal{B}}^{(c)}} \sum_{b \in \mathcal{B}} n_b^{(c)} AuthorInfluence_c(a \to b), & \text{if } N_{\mathcal{B}}^{(c)} > 0 \\ \tau_{a|\mathcal{V}}^{(c)}, & \text{otherwise,} \end{cases}$$

where $N_{\mathcal{B}}^{(c)}$ is the total number of occurrences of concept $c$ in the set $\mathcal{B}$, $n_b^{(c)}$ is the frequency of concept $c$ in paper $b$, and $\mathcal{V}$ is the set of all papers in the corpus. Here, the influence to each $b \in \mathcal{B}$ is weighted by the relative prevalence of concept $c$ with respect to $\mathcal{B}$, $n_b^{(c)}/N_{\mathcal{B}}^{(c)}$. We note that if a researcher's trusted set $\mathcal{B}$ contains no occurrences of a particular concept, we assign the trust value to $\tau_{a|\mathcal{V}}^{(c)}$, as if all the papers in the corpus were trusted equally.

In order to incorporate trust into paper selection, we assume an author will trust a paper if and only if he trusts *at least one of its authors*. This intuition can be formalized by defining a modified notion of set influence, where the researcher's preferences towards the authors are directly taken into account:

$$Influence_c(q \leftrightarrow \mathcal{A}|\mathcal{B}) = 1 - \prod_{d \in \mathcal{A}} \left( 1 - Influence_c(q \leftrightarrow d)\beta_d^{(c)} T_{d|\mathcal{B}}^{(c)} \right),$$

where $T_{d|\mathcal{B}}^{(c)} = 1 - \prod_{a \in authors(d)} (1 - \tau_{a|\mathcal{B}}^{(c)})$.

11

**Algorithm 2** Dynamic Programming Heuristic for Author Influence

---

$N$: number of documents
$\mathcal{C}$: vocabulary of concepts
// Initialize to empty 3D array
// $authorInfluence[c][a][y]$ will contain influence from author $a$ to paper $y$ w.r.t. concept $c$.
$authorInfluence \leftarrow array[|\mathcal{C}|][numAuthors][N]$
**for all** $c \in \mathcal{C}$ **do**
  **for all** authors $a$ **do**
    // Every author influences his or her own papers
    $authorInfluence[c][a][papers(a)] \leftarrow 1$
  $topoOrder \leftarrow$ topological order of nodes in $G_c$
  **for** $y \in topoOrder$ **do**
    // $authorInfluence[c][][x]$ already calculated for all $x \in parents(y)$
    **if** $parents(y) = \emptyset$ **then**
      continue
    $influenceFromParents \leftarrow array[|parents(y)|]$
    **for all** $x \in parents(y)$ **do**
      // Influence to the parent multiplied by the edge weight
      $influenceFromParents[x] \leftarrow authorInfluence[c][][x] \cdot \theta_{x \rightarrow y}^{(c)}$
    // Product is element-wise
    $authorInfluence[c][][y] \leftarrow 1 - \prod_{x \in parents(y)}(1 - influenceFromParents[x])$
    // Retain authors' self-influence
    $authorInfluence[c][authors(y)][y] \leftarrow 1$

---

We can now define our *personalized* objective function as:

$$F_{\mathcal{Q}|\mathcal{B}}(\mathcal{A}) \quad = \quad \sum_{q \in \mathcal{Q}} \sum_{c \in \mathcal{C}} \gamma_q^{(c)} Influence_c(q \leftrightarrow \mathcal{A}|\mathcal{B}). \tag{3}$$

Maximizing $F_{\mathcal{Q}|\mathcal{B}}(\mathcal{A})$ subject to $|\mathcal{A}| \leq k$, for some budget of $k$ papers, leads to a personalized set of papers tailored to someone who trusts $\mathcal{B}$. This function shares the same theoretical properties as Equation 2 and can be optimized efficiently in the same manner.

**Theorem 3.** *Equation 3 is submodular and monotonic.*

*Proof.* We use the exact same argument that we did for proving Equation 2 is submodular and monotonic, using the added fact that the document trust weight $T_{d|\mathcal{B}}^{(c)}$ is a probability in the range $[0, 1]$. $\qquad\square$

Figure 6 shows our PNAS example from before, with the same query set $\mathcal{Q} = \{4468, 5688\}$, but now incorporating the trust preferences of two hypothetical researchers: a plant biologist (A) and an immunologist (B). The figure highlights how differences in trust preferences can manifest themselves in article selection. In Figure 7, we provide another example, this time from computer science. Here, we take the famous Faloutsos, Faloutsos and Faloutsos paper, "On power-law relationships of the Internet topology" [19], and select related literature for it using the trust preferences of each author. Specifically, the visualization in the figure shows that by assuming that Michalis Faloutsos trusts SIGCOMM papers, Petros Faloutsos trusts SIGGRAPH papers, and Christos Faloutsos trusts KDD papers, we can select related work tailored to each author's perspective. While some relevant papers are common to all three points of view, other selected papers are particular to just one. For example, in Christos' data mining-focused result set, we find a few papers related to the evolution of social networks (e.g., "Microscopic evolution of social networks" by Leskovec et al.) which are not found in Michalis' and Petros' results. Moreover, these papers are not selected in the unpersonalized setting, when no trust preferences are taken into account.
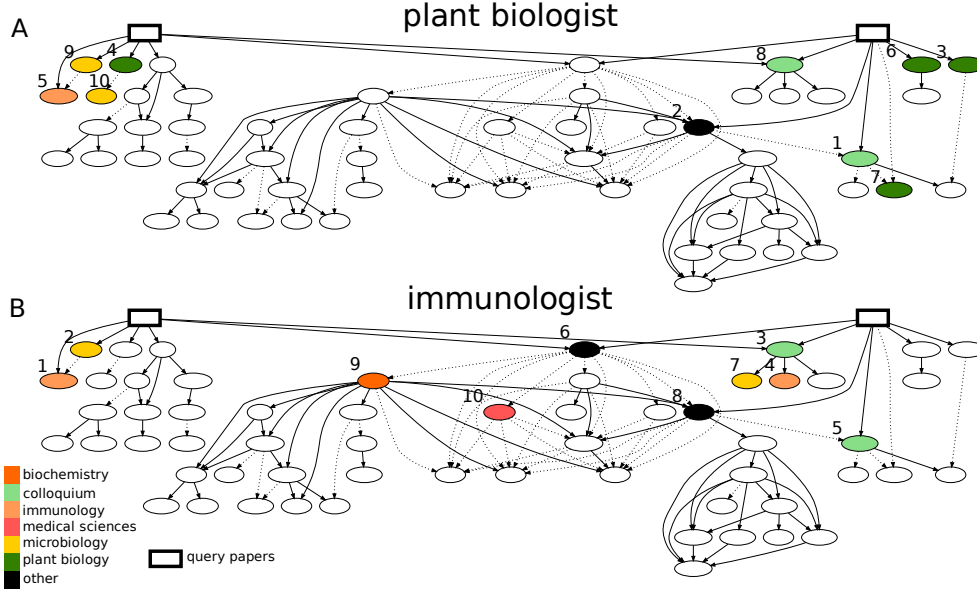
Figure 6: Top ten papers selected for $\mathcal{Q} = \{4468, 5688\}$ where $\mathcal{B}$ is defined as (A) all the plant biology papers, or (B) all the immunology papers, in the PNAS data set. Node colors correspond to article classification, as indicated by the key. (Colloquium refers to the National Academy of Sciences Colloquium on Virulence and Defense in Host-Pathogen Interactions: Common Features Between Plants and Animals. "Other" refers to unclassified papers, e.g., "From the Academy.".) Numbers indicate order of selection by optimization algorithm, roughly indicating order of importance (cf. Tables 3 and 4).

# 4 Approach Summary

We summarize our approach as follows:

**Initialization**

1. Define a vocabulary of concepts $\mathcal{C}$ (e.g., technical terms).
2. For each concept $c \in \mathcal{C}$, define a directed, acyclic graph $G_c$, with edge weights as in Section 2.1.
3. Compute relevance weights $\gamma_d^{(c)}$ and $\beta_d^{(c)}$, for all $c \in \mathcal{C}$ and documents $d$, as described in Section 2.3.1.
4. Precompute $Influence_c(u \leftrightarrow v)$ for all concepts $c$, and all pairs of documents $u$ and $v$, using Algorithm 1.
5. Similarly, precompute $AuthorInfluence_c(a \rightarrow b)$, for all authors $a$, all documents $b$, and all $c \in \mathcal{C}$, using Algorithm 2.

**Per user**

Given a user's trusted set of papers $\mathcal{B}$, compute $\tau_{a|\mathcal{B}}^{(c)}$ for all authors $a$ and $c \in \mathcal{C}$.

**Per query**

Given query set $\mathcal{Q}$, optimize Equation 3 using CELF [29].

# 5 Related Work

Researchers in both the library science and computer science communities have studied the shortcomings of the traditional keyword search paradigm [5, 35, 37]. In fact, our specific query model of defining a researcher's information need as a set of papers rather than as a keyword string has been described before [9, 30]. In one particularly related line of research, collaborative filtering techniques that have been successful for movie and product recommendations were adapted to the paper recommendation setting [30, 45]. Another approach uses

**no trust preferences**

**networks perspective**

**graphics perspective**

**data mining perspective**

Figure 7: A visualization of related work for Faloutsos, Faloutsos, and Faloutsos' "On power-law relationships of the Internet topology." The top word cloud represents papers selected using Equation 2, with no trust preferences. (The size of each word in the cloud is proportional to its prevalence in the selected documents.) The subsequent three word clouds represent papers selected using Equation 3 with three different trusted sets $\mathcal{B}$, one representing each author's perspective. Each word cloud visualizes the selected papers that are unique to each author's result set. For example, the bottom word cloud shows the papers found in Christos' data mining-focused results, but do not exist in Petros' or Michalis' result sets.

hypothesis testing to determine the articles that most influence each paper–the paper's *Information Genealogy*–based only on article text [43]. Unlike these previous approaches, our methodology is based on a *unified* model of text and citations that places special emphasis on the different trust preferences of individual researchers.

Previous work has also considered the more general, yet related, problem of taking positive examples of membership in a set and using them to expand the set [17, 22]. While such approaches have been applied to the domain of research literature, they do not explicitly model the particular characteristics of our problem, e.g., the effect of citations, publication venues and authorship.

Moreover, it is also important to note that our algorithm is *operational* in that it describes a method for selecting papers, in contrast with many *descriptive* studies in bibliometrics, sociology and other fields [13, 40, 33, 34, 4, 41]. In particular, the large body of work on *topic modeling* in computer science and statistics focuses on fitting probabilistic models to document collections by modeling latent themes in the data [8]. While often applied to
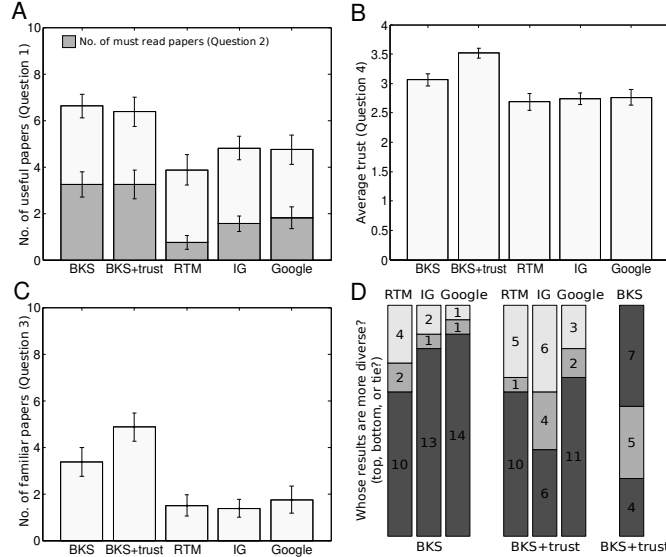
Figure 8: User study results comparing two variants of our algorithm, Beyond Keyword Search (BKS), with and without incorporating trust preferences, with the Relational Topic Model (RTM), Information Genealogy (IG) and Google Scholar. Values in bar plots (A), (B) and (C) are responses to the indicated study questions averaged over all sixteen participants, with error bars indicating one standard error. (D) shows how many participants (out of 16) found that our method produced more diverse results compared to the alternative techniques.

corpora of scholarly literature [18, 24, 3, 7, 42, 6, 15, 21], paper recommendation is not the primary objective of these models. Rather, our algorithm follows from a line of work that frames document selection as an explicit optimization problem (cf. [16]).

Finally, we note that the approach we describe in this paper is, in fact, agnostic to the specific definition of influence we use, and thus while we define influence to have an explicit probabilistic interpretation, other such definitions are possible. For instance, recent work by Lao and Cohen [28] provides an approach based on path-constrained random walks, which we can plug in as an alternative definition for influence.

# 6  Experimental Results

While these illustrative examples provide intuition, in order to truly evaluate our methodology we must solicit feedback from real scientific researchers. To this end, we conducted a user study involving sixteen subjects (all doctoral students in computer science or related fields).

We compare two variants of our algorithm–with and without incorporating trust preferences of the participant–with three representative alternative techniques: Google Scholar [23], Information Genealogy [43] (a hypothesis testing approach based on document text), and the Relational Topic Model [11] (a state-of-the-art topic model incorporating both text and citations to model latent themes in data)[3]. For each participant, we use each of these techniques to find related work for a previously written paper–that participant's *study paper*–thereby simulating a real research scenario. We define each query set $\mathcal{Q}$ to be the references of the corresponding study paper, and we ask each participant to list up to four trusted conferences or journals, which we use to define $\mathcal{B}$. The articles used in this study come from the ACM Digital Library [1], as described in the appendix.

---

[3]Previous work [30] has shown that Google keyword search outperforms collaborative filtering techniques for selecting useful papers, and thus we do not directly compare against these approaches.

In the case of Google Scholar, we ask a coauthor of the participant to provide the ideal keyword query he or she would use to find related work for the study paper. We enter this query into Google Scholar, and retrieve a result set containing the top ten papers that also appear in our ACM data set. In some cases, the keyword query provided was too specific, resulting in fewer than ten Google Scholar results.

For the Relational Topic Modeling approach, we fit the model to our data using the collapsed Gibbs sampling package provided by the authors [10]. We use K=50 topics, a burn-in of 750 samples, and collect our results over 750 additional samples. The parameters are set according to guidance from the first author (alpha=1/K, beta=4, eta=1/(size of vocabulary)). To select a set of related work, we compute the link probability from the study paper to each additional paper, and return the top ten most likely new links. We note that we give the model access to the abstract of the study paper–information that our algorithms do not have access to.

Finally, as the Information Genealogy model only takes into account document text, we provide the algorithm with the abstract of the study paper and retrieve the ten papers in the corpus with the most influence on the study paper. We use the same convex optimization package as used by the authors of the paper [31].

Unlike many previous studies, each participant was asked to evaluate all five comparison methods, rather than just a single technique. In total, 612 distinct papers were recommended using these five techniques across all sixteen participants.

Each participant was presented with the recommended articles for his or her study paper in a double-blind fashion, masking the identity of the technique used to select each paper. Participants were asked to answer questions on the usefulness, novelty and trustworthiness of each paper with respect to their research.[4] Additionally, participants were presented with entire result sets and asked to evaluate them in terms of diversity. Figure 8 shows the results of the study, from which we can glean the following main observations:

1. On average, users find the papers our algorithm selects to be more useful than those selected by the comparison techniques. The topic modeling approach performs especially poorly, with fewer than half of selected papers deemed useful.
2. Explicitly modeling the individual trust preferences of users leads to more trustworthy papers being selected. However, this comes at the expense of novelty in the selected articles, as researchers are more familiar with the work of authors they trust.
3. Our algorithm provides more diverse results than the comparison techniques, which is unsurprising, as our objective functions penalize redundancy.

# 7 Discussion

These results illustrate the success of our approach in recommending highly relevant literature personalized to the preferences of individual researchers, acting as a promising complement to keyword search. On a personal note, employing our approach during the writing of *this article* led us to related work from another subfield of computer science that we had not discovered using more traditional search methods [30, 45]. In closing, we believe that the challenges researchers face in expressing their information needs extend beyond scientific literature to domains like patents, law and news, and the work presented herein is a significant step towards addressing this general concern.

# 8 Acknowledgments

---

[4]Specific questions asked can be found in the appendix.

# References

[1] ACM Digital Library. `http://portal.acm.org`.

[2] R. Adler, J. Ewing, and P. Taylor. Citation statistics. *Statistical Science*, 24:1–14, 2009.

[3] E. M. Airoldi, E. A. Erosheva, S. E. Fienberg, C. Joutard, T. Love, and S. Shringarpure. Reconceptualizing the classification of PNAS articles. *Proceedings of the National Academy of Sciences USA*, 2010.

[4] A.-L. Barabási. On the topology of the scientific collaboration networks. *Physica A*, 311:590–614, 2002.

[5] M. Bates. The design of browsing and berrypicking techniques for the online search interface. *Online Review*, 13:407–424, 1989.

[6] D. M. Blei and J. Lafferty. Dynamic topic models. In *ICML*, 2006.

[7] D. M. Blei and J. Lafferty. A correlated topic model of *science*. *Ann. Appl. Stat.*, 1:17–35, 2007.

[8] D. M. Blei and J. Lafferty. *Topic Models*. Chapman and Hall, 2009.

[9] K. Bollacker, S. Lawrence, and C. L. Giles. Discovering relevant scientific literature on the Web. *IEEE Intelligent Systems and their Applications*, 15:42–47, 2000.

[10] J. Chang. Collapsed Gibbs sampling for topic models. `http://cran.r-project.org/web/packages/lda/`, 2010.

[11] J. Chang and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010.

[12] P. Chen, H. Xie, S. Maslov, and S. Redner. Finding scientific gems with Google. *Journal of Informetrics*, 1:8–15, 2007.

[13] D. J. de Solla Price. Networks of scientific papers. *Science*, 149:510:515, 1965.

[14] D. Diderot. In D. Diderot and J. d'Alembert, editors, *Encyclopedia, or a systematic dictionary of the sciences, arts and crafts*, Paris, 1755. Briasson, David, Le Breton, and Durand. (tr. from French).

[15] L. Dietz, S. Bickel, and T. Scheffer. Unsupervised prediction of citation influences. In *ICML*, 2007.

[16] K. El-Arini, G. Veda, D. Shahaf, and C. Guestrin. Turning down the noise in the blogosphere. In *KDD*, 2009.

[17] C. Elkan and K. Noto. Learning classifiers from only positive and unlabeled data. In *KDD*, 2008.

[18] E. A. Erosheva, S. E. Fienberg, and J. Lafferty. Mixed-membership models of scientific publications. *Proceedings of the National Academy of Sciences USA*, 101:5220–5227, 2004.

[19] M. Faloutsos, P. Faloutsos, and C. Faloutsos. On power-law relationships of the Internet topology. In *SIGCOMM*, 1999.

[20] E. Garfield. Citation analysis as a tool in journal evaluation. *Science*, 178:471–479, 1972.

[21] S. Gerrish and D. M. Blei. A language-based approach to measuring scholarly impact. In *ICML*, 2010.

[22] Z. Ghahramani and K. A. Heller. Bayesian sets. In *NIPS*, 2006.

[23] Google Scholar. `http://scholar.google.com`.

[24] T. Griffiths and M. Steyvers. Finding scientific topics. *Proceedings of the National Academy of Sciences USA*, 101:5228–5235, 2004.

[25] J. E. Hirsch. An index to quantify an individual's scientific research output. *Proceedings of the National Academy of Sciences USA*, 102:16569–16572, 2005.

[26] S. Khuller, A. Moss, and J. Naor. The budgeted maximum coverage problem. *Information Processing Letters*, 1999.

[27] J. Kleinberg. Authoritative sources in a hyperlinked environment. *Journal of the ACM*, 46:604–632, 1999.

[28] N. Lao and W. W. Cohen. Relational learning using a combination of path-constrained random walks. *Machine Learning*, 81(1):53–67, 2010.

[29] J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, and N. Glance. Cost-effective outbreak detection in networks. In *KDD*, 2007.

[30] S. M. McNee, I. Albert, D. Cosley, P. Gopalkrishnan, S. K. Lam, A. M. Rashid, J. A. Konstan, and J. Riedl. On the recommending of citations for research papers. In *CSCW*, 2002.

[31] Mosek. http://www.mosek.com.

[32] G. Nemhauser, L. Wolsey, and M. Fisher. An analysis of the approximations for maximizing submodular set functions. *Mathematical Programming*, 14:265–294, 1978.

[33] M. E. J. Newman. Scientific collaboration networks: I. network construction and fundamental results. *Phys. Rev. E*, 64:016131, 2001.

[34] M. E. J. Newman. The structure of scientific collaboration networks. *Proc. Natl. Acad. Sci. USA*, 98:404–409, 2001.

[35] C. Olston and E. H. Chi. Scenttrails: Integrating browsing and searching on the Web. *ACM Transactions on Computer-Human Interaction*, 10:177–197, 2003.

[36] L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank citation ranking: bringing order to the Web. Technical report, Stanford University InfoLab, 1999.

[37] S. Pandit and C. Olston. Navigation-aided retrieval. In *WWW*, 2007.

[38] J. S. Provan and M. O. Ball. The complexity of counting cuts and of computing the probability that a graph is connected. *SIAM J. Comput.*, 12:777–788, 1983.

[39] F. Radicchi, S. Fortunato, B. Markines, and A. Vespignani. Diffusion of scientific credits and the ranking of scientists. *Physical Review E*, 80:056103, 2009.

[40] S. Redner. How popular is your paper? an empirical study of the citation distribution. *Eur. Phys. J. B*, 4:131–134, 1998.

[41] M. Rosvall and C. T. Bergstrom. Maps of random walks on complex networks reveal community structure. *Proc. Natl. Acad. Sci. USA*, 105:1118–1123, 2008.

[42] M. Rozen-Zvi, T. Griffiths, M. Steyvers, and P. Smyth. The author-topic model for authors and documents. In *UAI*, 2004.

[43] B. Shaparenko and T. Joachims. Information genealogy: uncovering the flow of ideas in non-hyperlinked document databases. In *KDD*, 2007.

[44] Thomson Reuters Web of Knowledge. `http://wokinfo.com/about/whatitis`.

[45] R. Torres, S. M. McNee, M. Abel, J. A. Konstan, and J. Riedl. Enhancing digital libraries with TechLens+. In *JCDL*, 2004.

[46] L. Valiant. The complexity of enumeration and reliability problems. *SIAM J. Comput.*, 8:410–421, 1979.

# A    Data details and preprocessing

In this paper we refer to two data sets of scientific publications:

1. *PNAS Data*: Five years worth of articles from the Proceedings of the National Academy of Sciences (1997-2001; 13,648 papers).

2. *ACM Data*: A subset of the Association for Computing Machinery Digital Library, focused on papers in machine learning and related areas (1959-2009; 35,042 papers).

Both data sets include the title, authors, publication date, venue or publication track, citations, abstract and full text (when available) of each paper. The particular running example in our paper refers to the PNAS articles in Table 1.

Each data set was preprocessed to ensure the acyclicity of the citation graph, as well as to extract a vocabulary $\mathcal{C}$ of important concepts. The content of each paper is represented as a frequency vector of these selected concepts.

**Processing the citation graph.**    Based on simple chronology, one would expect a citation graph to be acyclic; after all, a researcher cannot cite a paper if it does not yet exist. However, this is not quite the case in practice. For instance, colleagues writing several papers simultaneously may cite each other, leading to doubly-connected pairs in the graph. As our algorithms rely on the acyclicity of the citation graph, we take the following steps to remove cycles:

1. Remove self cycles from the graph (i.e., edges that start and end at the same node).

2. Find the strongly connected components (SCCs) of the graph (i.e., maximal subgraphs such that for any two nodes $x, y$ in the subgraph, there is a path from $x$ to $y$ and a path from $y$ to $x$). In a directed acyclic graph (DAG), all the SCCs are of size one. However, this is generally not the case in real citation graphs.

3. For SCCs of size two (i.e., "I cite you and you cite me"), we employ the following heuristic to determine which edge to cut:

   - If the two papers were published in different years, have the later paper cite the earlier paper.
   - Else, if number of citations is different, have the lesser cited paper cite the more highly cited paper.
   - Else, pick one of the two edges uniformly at random.

4. While the previous step takes care of most cycles, a few peculiar cases with SCCs of size greater than two usually remain. There are few enough of these that we look at each such component individually, and manually decide which edges to cut.

Finally, recall that we augment the citation graph with edges indicating common authorship. In this step, we only connect papers that were written within five years of each other, as influence may tend to diminish over time. Moreover, when augmenting the graph with these edges, we ensure that we are not creating any cycles.

**Selecting concepts.** A typical corpus of scientific publications may contain tens of thousands of unique words, but only a fraction of them will be informative. Thus, working with the entire set of words rather than a particular subset can be wasteful. To this end, for each data set, we select a subset of words that we use as *concepts*:

- Ignore stop words (e.g., "the," "and," "of," etc.), words containing non-alphanumeric characters, and words that are too long ($> 20$ characters) or too short ($< 3$ characters).

- Of the remaining words, select the top 10,000 most frequent.

- Of these words, select ones that appear in at least 40 articles but fewer than 3,500 articles. If a word appears in too few articles, it is likely to be overly specific, while if it appears in a large fraction of articles, it is likely to be too general (e.g., the word "cell" for the case of PNAS, or "computer" for ACM). (These numbers are for the PNAS data set. For the larger ACM collection, we require words to appear in at least 100 documents and in no more than 8,000.)

- Finally, in an attempt to avoid selecting marginal words, we only select words such that when they appear in a document, they appear at least twice (on average).

In future work, more sophisticated concept selection can be investigated.

## B  User Study Details

The following questions were asked of each user study participant, for each article presented:
1. Assume you came across this paper while working on the study paper. From reading the title and abstract, would you have been inclined to:
   (a) continue reading the paper (even if just to skim), because you think it might be useful to the work of the study paper?
   (b) walk away (i.e., from the title and abstract alone, you can already tell that this paper is not useful to the work of the study paper)?
2. Do you feel that this paper would have been a *must read* for you when working on the study paper? (i.e., you would have read this paper carefully had you known about it, and perhaps would have cited it) [Yes, No]
3. Did you know about this paper before? [Yes, No]
4. Taking the authors and venue into account, would you be inclined to trust what this paper has to say?
   (a) For sure [4]
   (b) Probably [3]
   (c) Not sure [2]
   (d) Probably not [1]
   (e) Not at all [0]

(Figure 8A plots the responses to questions 1 and 2. Figure 8B plots the responses to question 4. Figure 8C plots the responses to question 3.)

After answering these questions for all papers selected by all five approaches, the participant is presented with all ten pairings of the five approaches, head to head, one pair of result sets at a time (e.g., RTM results on the left of the screen, our results with trust on the right). For each pair of result sets, the participant is asked to indicate which of the result sets is more diverse, or if they are equally diverse. As a diverse set of useless results is not beneficial to a researcher, in this part of the study we only display the papers that were indicated as useful by the participant in the previous section (i.e., an affirmative answer to question 1). (These diversity results are plotted in Figure 8D.)

## C  Selected Papers

Tables 2-4 show the papers selected for our running PNAS example from the main text. In particular, Tables 3 and 4 show the papers presented in Figure 6. Table 5 provides the papers selected for the example in Figure 1D.

Table 1: Articles from PNAS example

| ID | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 160 | Physiological reactions of nitric oxide and hemoglobin: A radical re-think | 1999 | 96 | 9967-9969 |
| 244 | Isolation of a temperate bacteriophage encoding the type III effector protein SopE from an epidemic *Salmonella typhimurium* strain | 1999 | 96 | 9845-9850 |
| 292 | Bacteriophages in the evolution of pathogen-host interactions | 1999 | 96 | 9452-9454 |
| 1139 | Nitroreductase A is regulated as a member of the *soxRS* regulon of *Escherichia coli* | 1999 | 96 | 3537-3539 |
| 1304 | A mechanism of paraquat toxicity involving nitric oxide synthase | 1999 | 96 | 12760-12765 |
| 1839 | Ancient origins of nitric oxide signaling in biological systems | 1999 | 96 | 14206-14207 |
| 2094 | Hemoglobin induction in mouse macrophages | 1999 | 96 | 6643-6647 |
| 2136 | Virulent *Salmonella typhimurium* has two periplasmic Cu, Zn-superoxide dismutases | 1999 | 96 | 7502-7507 |
| 2389 | A highly conserved sequence is a novel gene involved in *de novo* vitamin B6 biosynthesis | 1999 | 96 | 9374-9378 |
| 2452 | The oxyhemoglobin reaction of nitric oxide | 1999 | 96 | 9027-9032 |
| 4468 | Periplasmic superoxide dismutase protects *Salmonella* from products of phagocyte NADPH-oxidase and nitric oxide synthase | 1997 | 94 | 13997-14001 |
| 5550 | Nitric oxide in plant immunity | 1998 | 95 | 10345-10347 |
| 5688 | Defense gene induction in tobacco by nitric oxide, cyclic GMP, and cyclic ADP-ribose | 1998 | 95 | 10328-10333 |
| 7273 | Roles for mannitol and mannitol dehydrogenase in active oxygen-mediated plant defense | 1998 | 95 | 15129-15133 |
| 8305 | *S*-nitrosothiol repletion by an inhaled gas regulates pulmonary function | 2001 | 98 | 5792-5797 |
| 8365 | Flavohemoglobin denitrosylase catalyzes the reaction of a nitroxyl equivalent with molecular oxygen | 2001 | 98 | 10108-10112 |
| 8445 | Expression and phylogeny of claudins in vertebrate primordia | 2001 | 98 | 10196-10201 |
| 8490 | Peptide methionine sulfoxide reductase from *Escherichia coli* and *Mycobacterium tuberculosis* protects bacteria against oxidative damage from reactive nitrogen intermediates | 2001 | 98 | 9901-9906 |

Tables 6-9 show the papers selected for the example in Figure 7.

Table 1 (cont.): Articles from PNAS example

| ID | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 8643 | Plant mitogen-activated protein kinase cascades: Negative regulatory roles turn out positive | 2001 | 98 | 784-786 |
| 8853 | Myoglobin: A scavenger of bioactive NO | 2001 | 98 | 735-740 |
| 8901 | Simultaneous observation of the O—O and Fe—$O_2$ stretching modes in oxyhemoglobins | 2001 | 98 | 479-484 |
| 8910 | Activation of a mitogen-activated protein kinase pathway is involved in disease resistance in tobacco | 2001 | 98 | 741-746 |
| 9135 | Catalytic consumption of nitric oxide by 12/15- lipoxygenase: Inhibition of monocyte soluble guanylate cyclase activation | 2001 | 98 | 8006-8011 |
| 9318 | *Helicobacter pylori* arginase inhibits nitric oxide production by eukaryotic cells: A strategy for bacterial survival | 2001 | 98 | 13844-13849 |
| 9429 | Reciprocal electromechanical properties of rat prestin: The motor molecule from rat outer hair cells | 2001 | 98 | 4178-4183 |
| 9452 | B lymphocyte-restricted expression of prion protein does not enable prion replication in prion protein knockout mice | 2001 | 98 | 4034-4037 |
| 9467 | Plasma nitrite rather than nitrate reflects regional endothelial nitric oxide synthase activity but lacks intrinsic vasodilator action | 2001 | 98 | 12814-12819 |
| 9573 | Supermolecular structure of the enteropathogenic *Escherichia coli* type III secretion system and its direct interaction with the EspA-sheath-like structure | 2001 | 98 | 11638-11643 |
| 9582 | Modulation of nitric oxide bioavailability by erythrocytes | 2001 | 98 | 11771-11776 |
| 9625 | Cysteine-3635 is responsible for skeletal muscle ryanodine receptor modulation by NO | 2001 | 98 | 11158-11162 |
| 9890 | *In vivo* mechanism-based inactivation of *S*-adenosylmethionine decarboxylases from *Escherichia coli*, *Salmonella typhimurium*, and *Saccharomyces cerevisiae* | 2001 | 98 | 10578-10583 |
| 10008 | Structure of sortase, the transpeptidase that anchors proteins to the cell wall of *Staphylococcus aureus* | 2001 | 98 | 6056-6061 |
| 10090 | Comparison of a hair bundle's spontaneous oscillations with its response to mechanical stimulation reveals the underlying active process | 2001 | 98 | 14380-14385 |
| 10118 | Compressive nonlinearity in the hair bundle's active response to mechanical stimulation | 2001 | 98 | 14386-14391 |
| 10123 | *In vivo* evidence for a cochlear amplifier in the hair-cell bundle of lizards | 2001 | 98 | 2826-2831 |

Table 1 (cont.): Articles from PNAS example

| ID | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 10161 | Defective localization of the NADPH phagocyte oxidase to *Salmonella*-containing phagosomes in tumor necrosis factor p55 receptor-deficient macrophages | 2001 | 98 | 2561-2565 |
| 10372 | Regulation of the *Mycobacterium tuberculosis* hypoxic response gene encoding $\alpha$-crystallin | 2001 | 98 | 7534-7539 |
| 10605 | Physical basis of two-tone interference in hearing | 2001 | 98 | 9080-9085 |
| 10642 | A fatty acid desaturase modulates the activation of defense signaling pathways in plants | 2001 | 98 | 9448-9453 |
| 10693 | Scrapie prion protein accumulation by scrapie-infected neuroblastoma cells abrogated by exposure to a prion protein antibody | 2001 | 98 | 9295-9299 |
| 10844 | Neuroglobin is up-regulated by and protects neurons from hypoxic-ischemic injury | 2001 | 98 | 15306-15311 |
| 10850 | Oxygen radical inhibition of nitric oxide-dependent vascular function in sickle cell disease | 2001 | 98 | 15215-15220 |
| 10900 | Epitope tagging of chromosomal genes in *Salmonella* | 2001 | 98 | 15264-15269 |
| 10940 | Polymerization of a single protein of the pathogen *Yersinia enterocolitica* into needles punctures eukaryotic cells | 2001 | 98 | 4669-4674 |
| 11134 | Relative role of heme nitrosylation and $\beta$-cysteine 93 nitrosation in the transport and metabolism of nitric oxide by hemoglobin in the human circulation | 2000 | 97 | 9943-9948 |
| 11770 | Protection from nitrosative stress by yeast flavohemoglobin | 2000 | 97 | 4672-4676 |
| 11791 | The *Pseudomonas syringae* Hrp pathogenicity island has a tripartite mosaic structure composed of a cluster of type III secretion genes bounded by exchangeable effector and conserved effector loci that contribute to parasitic fitness and pathogenicity in plants | 2000 | 97 | 4856-4861 |
| 12134 | *Arabidopsis* RelA/SpoT homologs implicate (p)ppGpp in plant signaling | 2000 | 97 | 3747-3752 |
| 12176 | Cochlear mechanisms from a phylogenetic viewpoint | 2000 | 97 | 11736-11743 |
| 12270 | Putting ion channels to work: Mechanoelectrical transduction, adaptation, and amplification by hair cells | 2000 | 97 | 11765-11772 |
| 12286 | Molecular mechanisms of sound amplification in the mammalian cochlea | 2000 | 97 | 11759-11764 |
| 12379 | Contribution of *Salmonella typhimurium* type III secretion components to needle complex formation | 2000 | 97 | 11008-11013 |
| 13042 | A conserved amino acid sequence directing intracellular type III secretion by *Salmonella typhimurium* | 2000 | 97 | 7539-7544 |

Table 1 (cont.): Articles from PNAS example

| ID | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 13204 | Reactive oxygen and nitrogen intermediates in the relationship between mammalian hosts and microbial pathogens | 2000 | 97 | 8841-8848 |
| 13240 | The *Arabidopsis dnd1* "defense, no death" gene encodes a mutated cyclic nucleotide-gated ion channel | 2000 | 97 | 9323-9328 |
| 13264 | Nitric oxide and salicylic acid signaling in plant defense | 2000 | 97 | 8849-8855 |
| 13279 | Genetic complexity of pathogen perception by plants: The example of *Rcr3*, a tomato gene required specifically by *Cf-2* | 2000 | 97 | 8807-8814 |
| 13283 | *Pseudomonas syringae* Hrp type III secretion system and effector proteins | 2000 | 97 | 8770-8777 |
| 13316 | Nitric oxide prevents cardiovascular disease and determines survival in polyglobulic mice overexpressing erythropoietin | 2000 | 97 | 11609-11613 |
| 13344 | Role of circulating nitrite and *S*-nitrosohemoglobin in the regulation of regional blood flow in humans | 2000 | 97 | 11482-11487 |

Table 2: Selected papers for PNAS example (no trust)

| Rank | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 1 | Nitric oxide in plant immunity | 1998 | 95 | 10345-10347 |
| 2 | Defective localization of the NADPH phagocyte oxidase to *Salmonella*-containing phagosomes in tumor necrosis factor p55 receptor-deficient macrophages | 2001 | 98 | 2561-2565 |
| 3 | Ancient origins of nitric oxide signaling in biological systems | 1999 | 96 | 14206-14207 |
| 4 | Virulent *Salmonella typhimurium* has two periplasmic Cu, Zn-superoxide dismutases | 1999 | 96 | 7502-7507 |
| 5 | Nitric oxide and salicylic acid signaling in plant defense | 2000 | 97 | 8849-8855 |
| 6 | Reactive oxygen and nitrogen intermediates in the relationship between mammalian hosts and microbial pathogens | 2000 | 97 | 8841-8848 |
| 7 | A mechanism of paraquat toxicity involving nitric oxide synthase | 1999 | 96 | 12760-12765 |
| 8 | Roles for mannitol and mannitol dehydrogenase in active oxygen-mediated plant defense | 1998 | 95 | 15129-15133 |
| 9 | The *Arabidopsis dnd1* "defense, no death" gene encodes a mutated cyclic nucleotide-gated ion channel | 2000 | 97 | 9323-9328 |
| 10 | *Arabidopsis* RelA/SpoT homologs implicate (p)ppGpp in plant signaling | 2000 | 97 | 3747-3752 |

Table 3: Selected papers for PNAS example (as a plant biologist)

| Rank | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 1 | Nitric oxide and salicylic acid signaling in plant defense | 2000 | 97 | 8849-8855 |
| 2 | Ancient origins of nitric oxide signaling in biological systems | 1999 | 96 | 14206-14207 |
| 3 | The *Arabidopsis dnd1* "defense, no death" gene encodes a mutated cyclic nucleotide-gated ion channel | 2000 | 97 | 9323-9328 |
| 4 | Roles for mannitol and mannitol dehydrogenase in active oxygen-mediated plant defense | 1998 | 95 | 15129-15133 |
| 5 | Defective localization of the NADPH phagocyte oxidase to *Salmonella*-containing phagosomes in tumor necrosis factor p55 receptor-deficient macrophages | 2001 | 98 | 2561-2565 |
| 6 | *Arabidopsis* RelA/SpoT homologs implicate (p)ppGpp in plant signaling | 2000 | 97 | 3747-3752 |
| 7 | A fatty acid desaturase modulates the activation of defense signaling pathways in plants | 2001 | 98 | 9448-9453 |
| 8 | Reactive oxygen and nitrogen intermediates in the relationship between mammalian hosts and microbial pathogens | 2000 | 97 | 8841-8848 |
| 9 | Virulent *Salmonella typhimurium* has two periplasmic Cu, Zn-superoxide dismutases | 1999 | 96 | 7502-7507 |
| 10 | A highly conserved sequence is a novel gene involved in *de novo* vitamin B6 biosynthesis | 1999 | 96 | 9374-9378 |

Table 4: Selected papers for PNAS example (as an immunologist)

| Rank | Title | Year | Volume | Pages |
|---|---|---|---|---|
| 1 | Defective localization of the NADPH phagocyte oxidase to *Salmonella*-containing phagosomes in tumor necrosis factor p55 receptor-deficient macrophages | 2001 | 98 | 2561-2565 |
| 2 | Virulent *Salmonella typhimurium* has two periplasmic Cu, Zn-superoxide dismutases | 1999 | 96 | 7502-7507 |
| 3 | Reactive oxygen and nitrogen intermediates in the relationship between mammalian hosts and microbial pathogens | 2000 | 97 | 8841-8848 |
| 4 | *Helicobacter pylori* arginase inhibits nitric oxide production by eukaryotic cells: A strategy for bacterial survival | 2001 | 98 | 13844-13849 |
| 5 | Nitric oxide and salicylic acid signaling in plant defense | 2000 | 97 | 8849-8855 |
| 6 | Nitric oxide in plant immunity | 1998 | 95 | 10345-10347 |
| 7 | Peptide methionine sulfoxide reductase from *Escherichia coli* and *Mycobacterium tuberculosis* protects bacteria against oxidative damage from reactive nitrogen intermediates | 2001 | 98 | 9901-9906 |
| 8 | Ancient origins of nitric oxide signaling in biological systems | 1999 | 96 | 14206-14207 |
| 9 | The oxyhemoglobin reaction of nitric oxide | 1999 | 96 | 9027-9032 |
| 10 | A mechanism of paraquat toxicity involving nitric oxide synthase | 1999 | 96 | 12760-12765 |

Table 5: Selected papers for example in Figure 1D

| Rank | Title | Year | Volume | Pages |
|------|-------|------|--------|-------|
| 1 | Defense gene induction in tobacco by nitric oxide, cyclic GMP, and cyclic ADP-ribose | 1998 | 95 | 10328-10333 |
| 2 | Ancient origins of nitric oxide signaling in biological systems | 1999 | 96 | 14206-14207 |
| 3 | Periplasmic superoxide dismutase protects *Salmonella* from products of phagocyte NADPH-oxidase and nitric oxide synthase | 1997 | 94 | 13997-14001 |
| 4 | A mechanism of paraquat toxicity involving nitric oxide synthase | 1999 | 96 | 12760-12765 |
| 5 | Nitroreductase A is regulated as a member of the *soxRS* regulon of *Escherichia coli* | 1999 | 96 | 3537-3539 |
| 6 | Nitric oxide and salicylic acid signaling in plant defense | 2000 | 97 | 8849-8855 |
| 7 | *S*-nitrosothiol repletion by an inhaled gas regulates pulmonary function | 2001 | 98 | 5792-5797 |
| 8 | Cysteine-3635 is responsible for skeletal muscle ryanodine receptor modulation by NO | 2001 | 98 | 11158-11162 |
| 9 | The oxyhemoglobin reaction of nitric oxide | 1999 | 96 | 9027-9032 |
| 10 | Protection from nitrosative stress by yeast flavohemoglobin | 2000 | 97 | 4672-4676 |
| 11 | Hemoglobin induction in mouse macrophages | 1999 | 96 | 6643-6647 |
| 12 | Physiological reactions of nitric oxide and hemoglobin: A radical rethink | 1999 | 96 | 9967-9969 |
| 13 | Cochlear mechanisms from a phylogenetic viewpoint | 2000 | 97 | 11736-11743 |
| 14 | Plant mitogen-activated protein kinase cascades: Negative regulatory roles turn out positive | 2001 | 98 | 784-786 |
| 15 | Flavohemoglobin denitrosylase catalyzes the reaction of a nitroxyl equivalent with molecular oxygen | 2001 | 98 | 10108-10112 |
| 16 | Relative role of heme nitrosylation and $\beta$-cysteine 93 nitrosation in the transport and metabolism of nitric oxide by hemoglobin in the human circulation | 2000 | 97 | 9943-9948 |
| 17 | Role of circulating nitrite and *S*-nitrosohemoglobin in the regulation of regional blood flow in humans | 2000 | 97 | 11482-11487 |
| 18 | Modulation of nitric oxide bioavailability by erythrocytes | 2001 | 98 | 11771-11776 |
| 19 | Nitric oxide prevents cardiovascular disease and determines survival in polyglobulic mice overexpressing erythropoietin | 2000 | 97 | 11609-11613 |
| 20 | Plasma nitrite rather than nitrate reflects regional endothelial nitric oxide synthase activity but lacks intrinsic vasodilator action | 2001 | 98 | 12814-12819 |

Table 6: Selected papers for example in Figure 7 (unpersonalized)

| Rank | Title | Authors | Year |
|---|---|---|---|
| 1 | Prediction of future world wide web traffic characteristics for capacity planning | Christensen, Javagal | 1997 |
| 2 | Self-similarity in World Wide Web traffic: evidence and possible causes | Crovella, Bestavros | 1997 |
| 3 | Characteristics of WWW Client-based Traces | Cunha et al. | 1995 |
| 4 | Empirically derived analytic models of wide-area TCP connections | Paxson | 1994 |
| 5 | End-to-end available bandwidth as a random autocorrelated QoS-relevant time-series | Chobanyan et al. | 2008 |
| 6 | Efficiently serving dynamic data at highly accessed web sites | Challenger et al. | 2004 |
| 7 | A Prefetching Protocol Using Client Speculation for the WWW | Bestavros, Cunha | 1995 |
| 8 | Power laws and the AS-level internet topology | Siganos et al. | 2003 |
| 9 | Power-law relationship and self-similarity in the itemset support distribution: analysis and applications | Chuang et al. | 2008 |
| 10 | On the origin of power laws in Internet topologies | Medina et al. | 2000 |
| 11 | Spatio-temporal network anomaly detection by assessing deviations of empirical measures | Paschalidis, Smaragdakis | 2009 |
| 12 | Network topology generators: degree-based vs. structural | Tangmunarunkit et al. | 2002 |
| 13 | Mathematical models for academic webs: linear relationship or non-linear power law? | Payne, Thelwall | 2005 |
| 14 | A random graph model for massive graphs | Aiello et al. | 2000 |

Table 7: Selected papers for example in Figure 7 (networks)

| Rank | Title | Authors | Year |
|---|---|---|---|
| 1 | Self-similarity in World Wide Web traffic: evidence and possible causes | Crovella, Bestavros | 1997 |
| 2 | Empirically derived analytic models of wide-area TCP connections | Paxson | 1994 |
| 3 | Power laws and the AS-level internet topology | Siganos et al. | 2003 |
| 4 | Characteristics of WWW Client-based Traces | Cunha et al. | 1995 |
| 5 | Weighted graphs and disconnected components: patterns and a generator | McGlohon et al. | 2008 |
| 6 | Learning for accurate classification of real-time traffic | Li, Moore | 2006 |
| 7 | BLINC: multilevel traffic classification in the dark | Karagiannis et al. | 2005 |
| 8 | Graphs over time: densification laws, shrinking diameters and possible explanations | Leskovec et al. | 2005 |
| 9 | Graph evolution: Densification and shrinking diameters | Leskovec et al. | 2007 |
| 10 | A Prefetching Protocol Using Client Speculation for the WWW | Bestavros, Cunha | 1995 |
| 11 | Scalable modeling of real graphs using Kronecker multiplication | Leskovec, Faloutsos | 2007 |
| 12 | ANF: a fast and scalable tool for data mining in massive graphs | Palmer et al. | 2002 |
| 13 | A random graph model for massive graphs | Aiello et al. | 2000 |
| 14 | Profiling internet backbone traffic: behavior models and applications | Xu et al. | 2005 |

Table 8: Selected papers for example in Figure 7 (graphics)

| Rank | Title | Authors | Year |
|------|-------|---------|------|
| 1 | Characteristics of WWW Client-based Traces | Cunha et al. | 1995 |
| 2 | Power laws and the AS-level internet topology | Siganos et al. | 2003 |
| 3 | Empirically derived analytic models of wide-area TCP connections | Paxson | 1994 |
| 4 | ANF: a fast and scalable tool for data mining in massive graphs | Palmer et al. | 2002 |
| 5 | Self-similarity in World Wide Web traffic: evidence and possible causes | Crovella, Bestavros | 1997 |
| 6 | Weighted graphs and disconnected components: patterns and a generator | McGlohon et al. | 2008 |
| 7 | Parallax photography: creating 3D cinematic effects from stills | Zheng et al. | 2009 |
| 8 | On inferring autonomous system relationships in the internet | Gao | 2001 |
| 9 | Power-law relationship and self-similarity in the itemset support distribution: analysis and applications | Chuang et al. | 2008 |
| 10 | On the origin of power laws in Internet topologies | Medina et al. | 2000 |
| 11 | Composable controllers for physics-based character animation | Faloutsos et al. | 2001 |
| 12 | Segmenting motion capture data into distinct behaviors | Barbič et al. | 2004 |
| 13 | Efficiently serving dynamic data at highly accessed web sites | Challenger et al. | 2004 |
| 14 | Graph mining: Laws, generators, and algorithms | Chakrabarti, Faloutsos | 2006 |

Table 9: Selected papers for example in Figure 7 (data mining)

| Rank | Title | Authors | Year |
|---|---|---|---|
| 1 | Characteristics of WWW Client-based Traces | Cunha et al. | 1995 |
| 2 | Power laws and the AS-level internet topology | Siganos et al. | 2003 |
| 3 | Graph evolution: Densification and shrinking diameters | Leskovec et al. | 2007 |
| 4 | Graphs over time: densification laws, shrinking diameters and possible explanations | Leskovec et al. | 2005 |
| 5 | Weighted graphs and disconnected components: patterns and a generator | McGlohon et al. | 2008 |
| 6 | Self-similarity in World Wide Web traffic: evidence and possible causes | Crovella, Bestavros | 1997 |
| 7 | Microscopic evolution of social networks | Leskovec et al. | 2008 |
| 8 | Statistical properties of community structure in large social and information networks | Leskovec et al. | 2008 |
| 9 | Scalable modeling of real graphs using Kronecker multiplication | Leskovec, Faloutsos | 2007 |
| 10 | Empirically derived analytic models of wide-area TCP connections | Paxson | 1994 |
| 11 | ANF: a fast and scalable tool for data mining in massive graphs | Palmer et al. | 2002 |
| 12 | Structure and evolution of online social networks | Kumar et al. | 2006 |
| 13 | Visualization of large networks with min-cut plots, A-plots and R-MAT | Chakrabarti et al. | 2007 |
| 14 | GraphScope: parameter-free mining of large time-evolving graphs | Sun et al. | 2007 |

**ML**

**MACHINE LEARNING**
**D E P A R T M E N T**

Carnegie Mellon University
5000 Forbes Avenue
Pittsburgh, PA 15213

**Carnegie Mellon**®